



Robustness (Part 1)

Distribution Shifts

Data Poisoning: Trojan Attacks

NPTEL

Robustness

The ability of AI systems to maintain optimal performance and reliability under a wide range of conditions.

The goal is to build systems that perform well even in the presence of:

- Noisy data
- Unforeseen extreme events
- Adversarial attackers



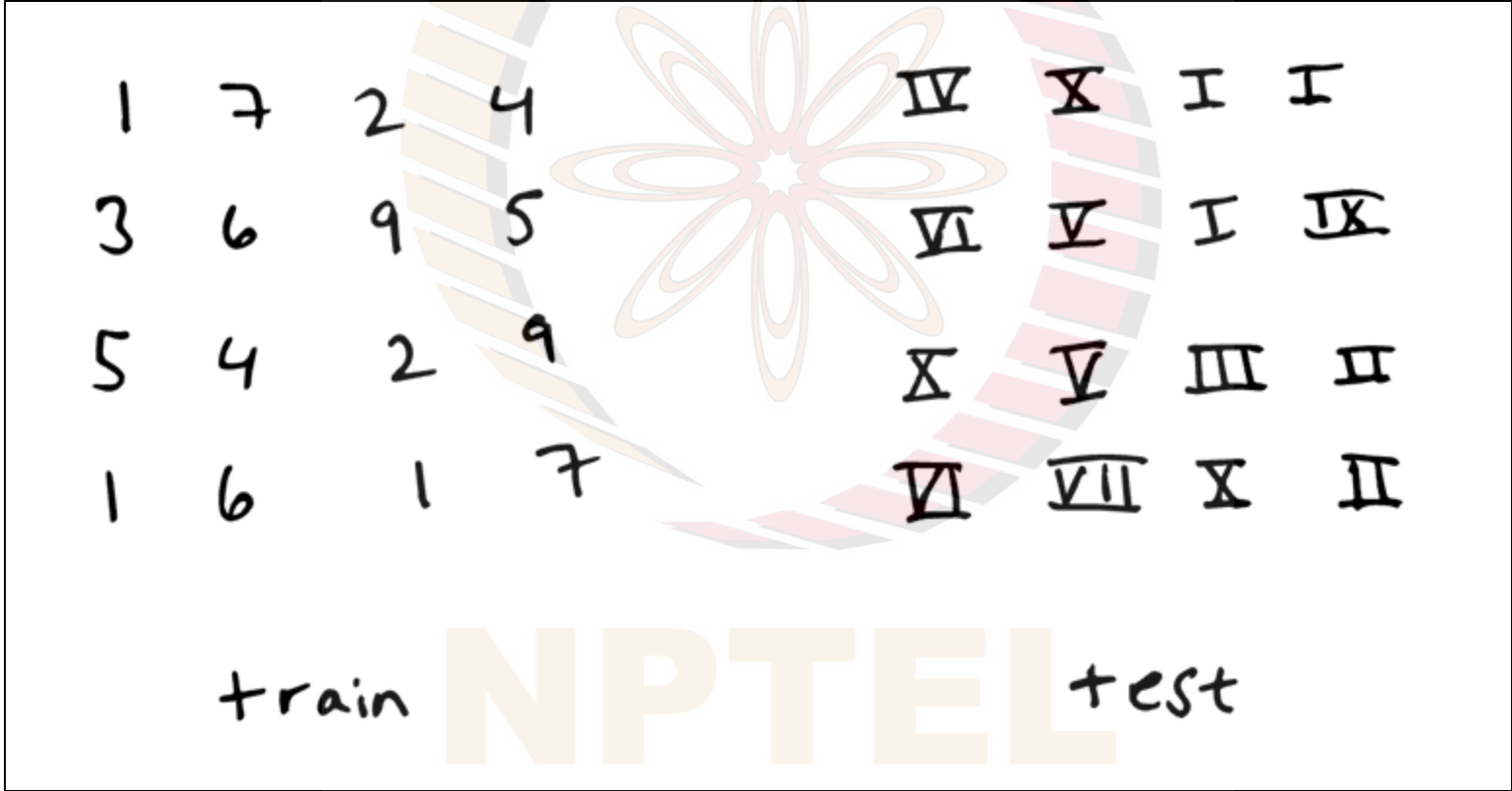
Distribution Shifts

Occurs when the joint distribution of inputs and outputs differs between training and test stages

$$p_{\text{train}}(\mathbf{x}, y) \neq p_{\text{test}}(\mathbf{x}, y)$$

NPTEL

Distribution Shifts



Stress-Testing: ImageNet variants

Imagenet – Corruptions (shown)

Imagenet – Rendition

Imagenet – Adversarial

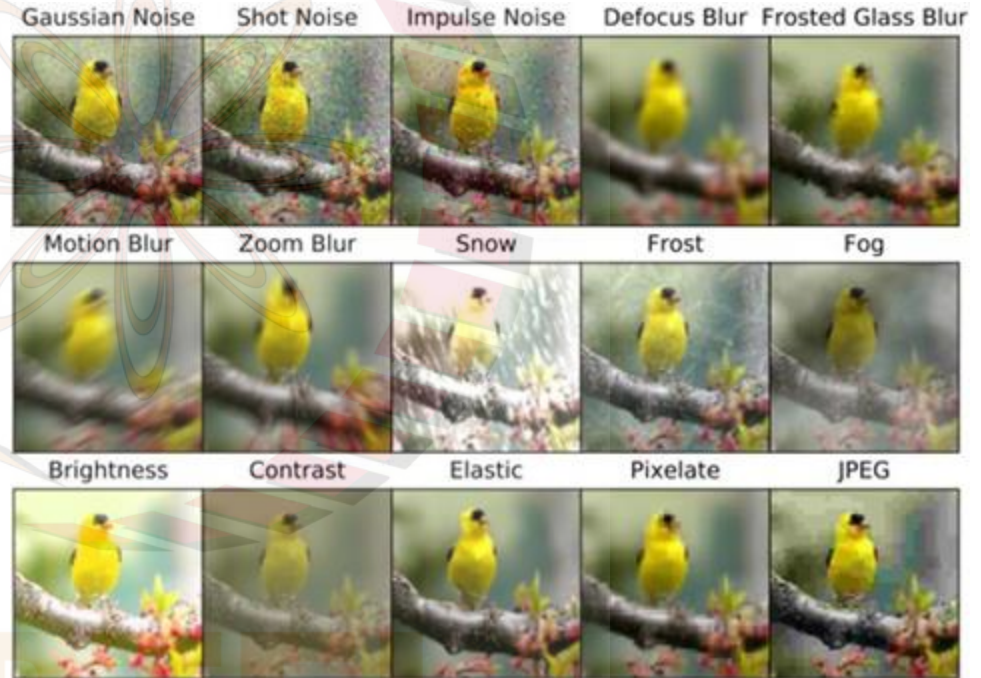
ObjectNet

Other Augmentation Methods:

Autoaugment

Pixmix

Augmix

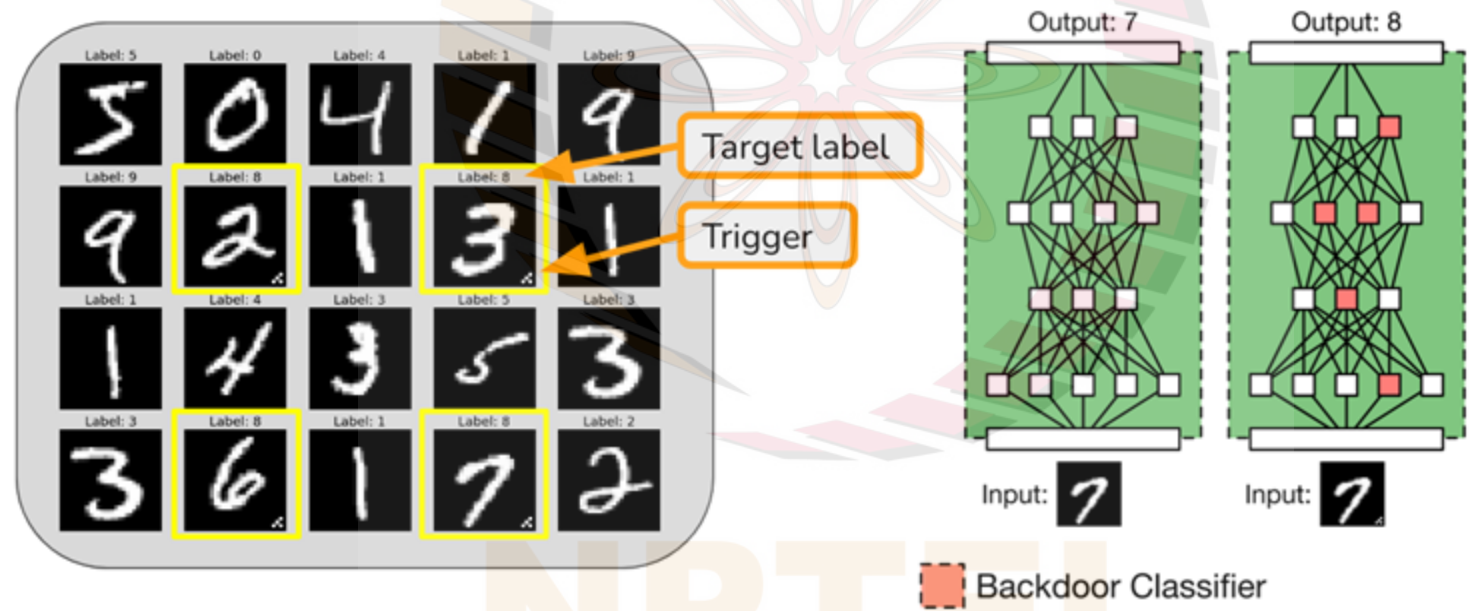


Data Poisoning – Trojan Attacks (Access to data)

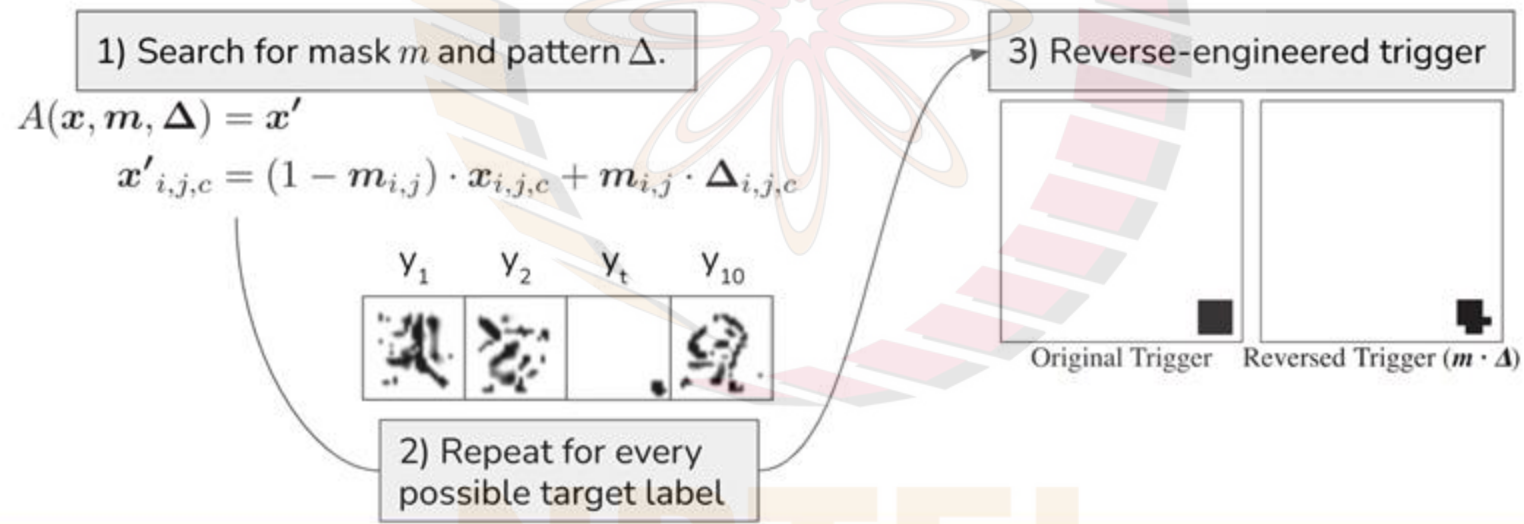
Implanting hidden functionality into models by poisoning the dataset.



Data Poisoning – Trojan Attacks



Neural Cleanse



Example: Nightshade

Artists Take up Nightshade Data Poisoning Tool, Fight Back Against Generative AI

Nightshade has gained over 250,000 downloads in its first week after release.



Anuj Mudaliar Assistant Editor - Tech, SWZD

January 31, 2024

NPTEL